

# **DiGA: Distil to Generalize and then Adapt for Domain Adaptive Semantic Segmentation**

Fengyi Shen<sup>1,2,3</sup>, Akhil Gurram<sup>2</sup>, Ziyuan Liu<sup>2</sup>, He Wang<sup>3,\*</sup>, Alois C. Knoll<sup>1,\*</sup>

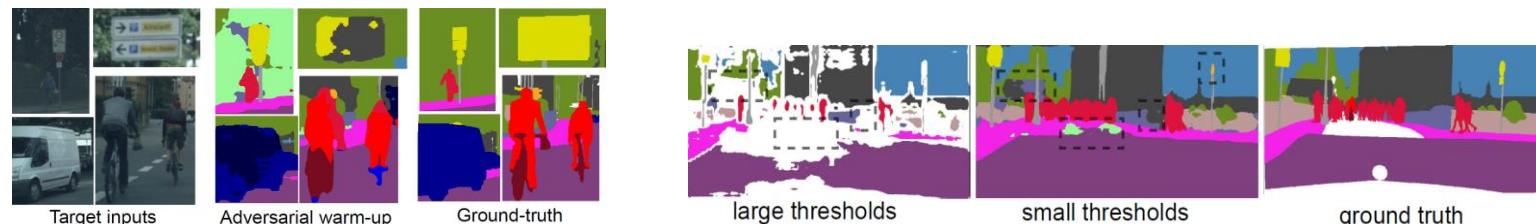
(<sup>1</sup>Technische Universität München, <sup>2</sup>Huawei Munich Research Center, <sup>3</sup> EPIC Lab, Peking University)

\*corresponding author

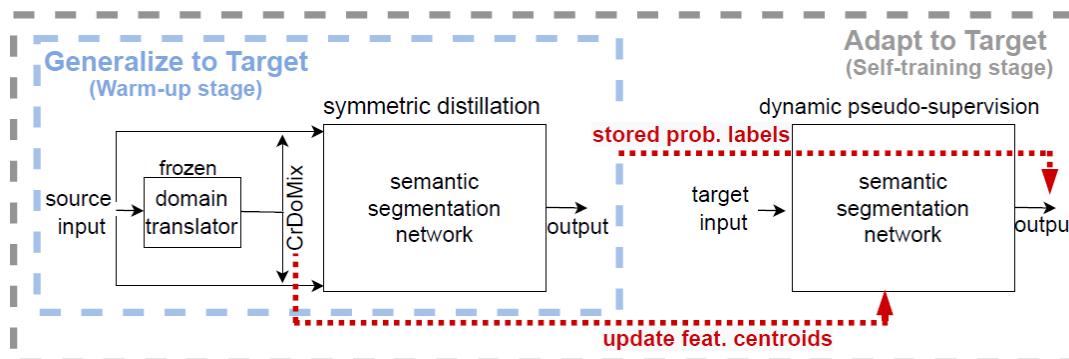
**CVPR'23 Poster WED-PM-335**

<https://github.com/fy-vision/DiGA>

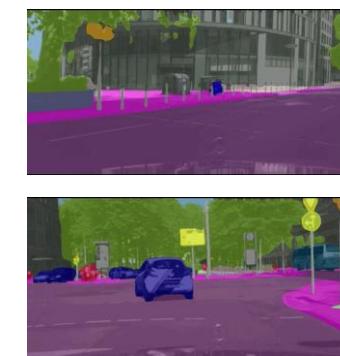
# Overview



## Stage-wise Training



DiGA Overview



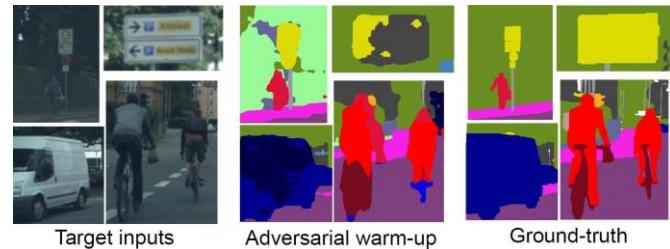
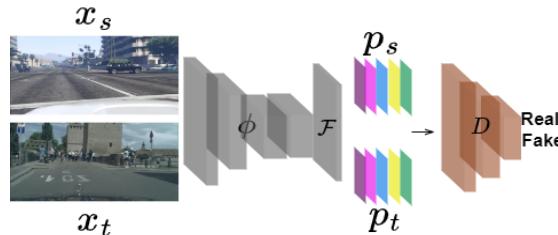
DiGA Demo

# Motivation

## Stage-wise training

### Warm-up:

AdaptSegNet (CVPR18), BDL (CVPR19)  
IAST (ECCV20), ProDA (CVPR21), CPSL (CVPR22) .....

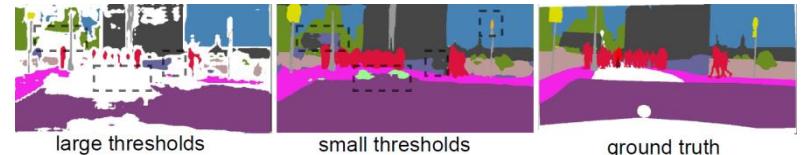
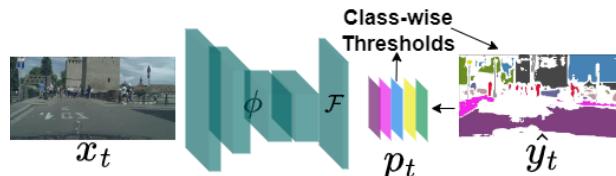


- *Class-unaware X*
- *Blind alignment X*

**Question1:** Can the model be trained w/o blindly aligning the target and source features in the warm-up stage?

### Self-training (pseudo-label):

BDL (CVPR19), CAG (NeurIPS19), IAST (ECCV20), ProDA (CVPR21)  
CPSL (CVPR22), BCL (ECCV22), ProCA (ECCV22) .....

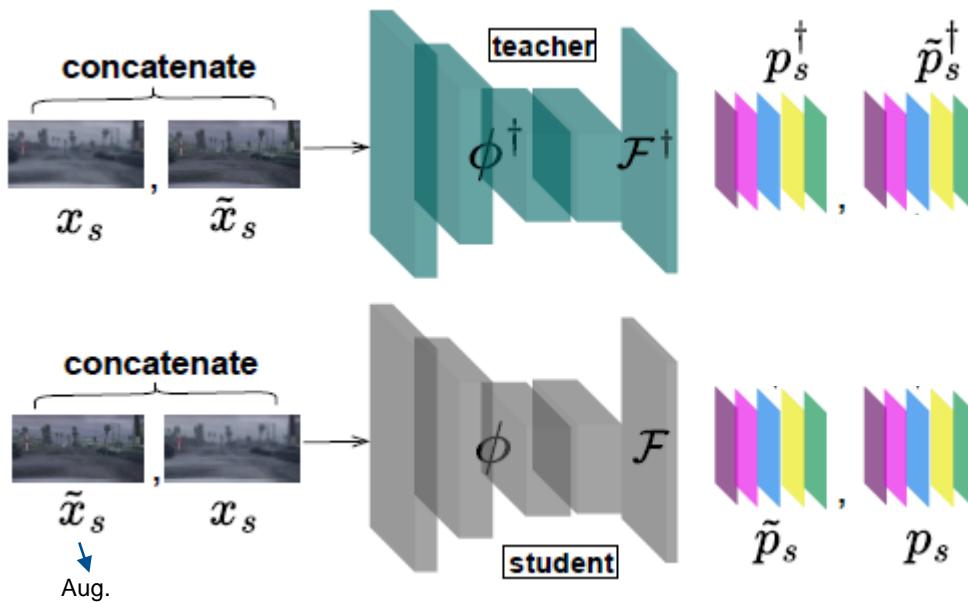


- *Large: insufficient learning X*
- *Small: introduce noise X*

**Question2:** Is it possible to avoid looking for thresholds during the self-training stage?

# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:



$$\{p_s^\dagger, \tilde{p}_s^\dagger\} = \sigma(\mathcal{F}^\dagger \circ \phi^\dagger(\{x_s, \tilde{x}_s\})) \quad (2)$$

$$\{\tilde{p}_s, p_s\} = \sigma(\mathcal{F} \circ \phi(\{\tilde{x}_s, x_s\})) \quad (3)$$

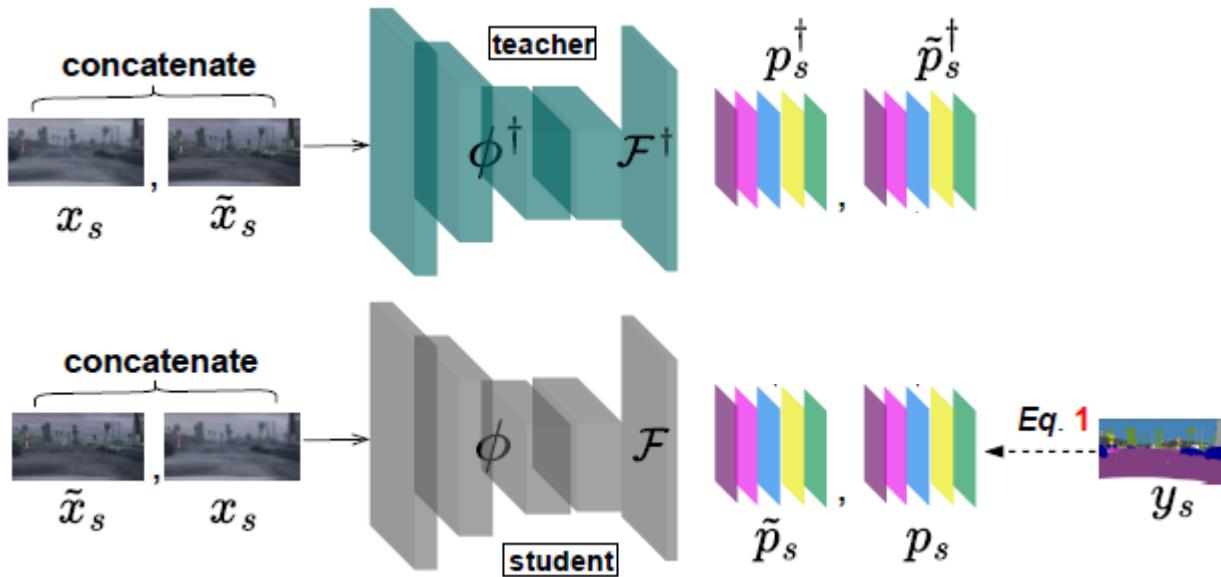
# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:

$$\mathcal{L}_s^{seg} = \sum_{h,w} \sum_c -y_s^{(c,h,w)} \log(p_s)^{(c,h,w)} \quad (1)$$

$$\{p_s^\dagger, \tilde{p}_s^\dagger\} = \sigma(\mathcal{F}^\dagger \circ \phi^\dagger(\{x_s, \tilde{x}_s\})) \quad (2)$$

$$\{\tilde{p}_s, p_s\} = \sigma(\mathcal{F} \circ \phi(\{\tilde{x}_s, x_s\})) \quad (3)$$



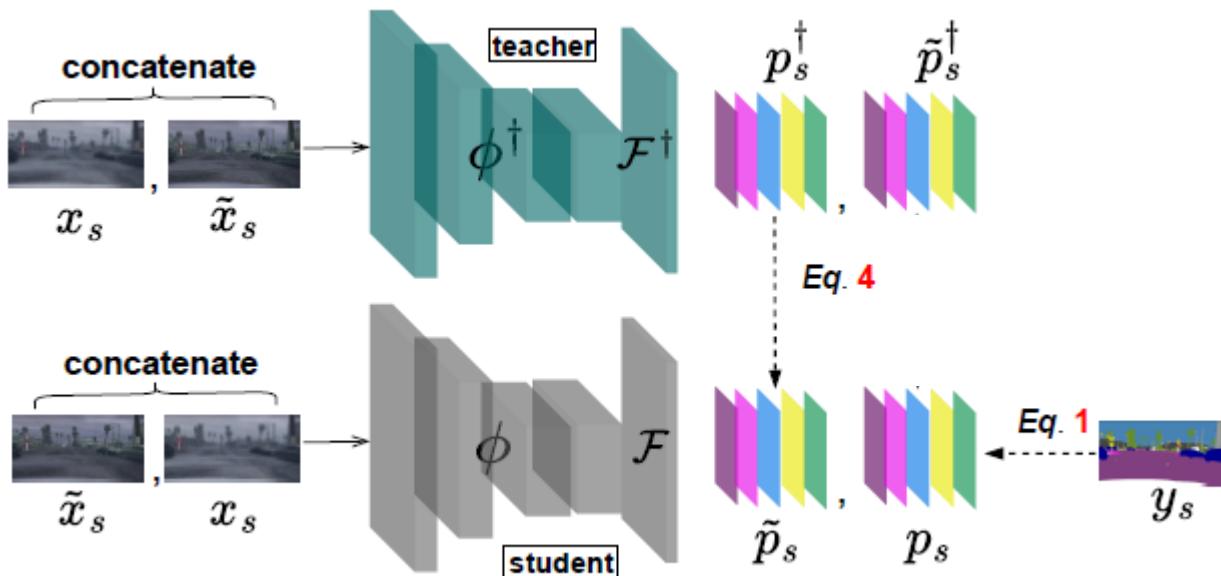
# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:

$$\mathcal{L}_s^{seg} = \sum_{h,w} \sum_c -y_s^{(c,h,w)} \log(p_s)^{(c,h,w)} \quad (1)$$

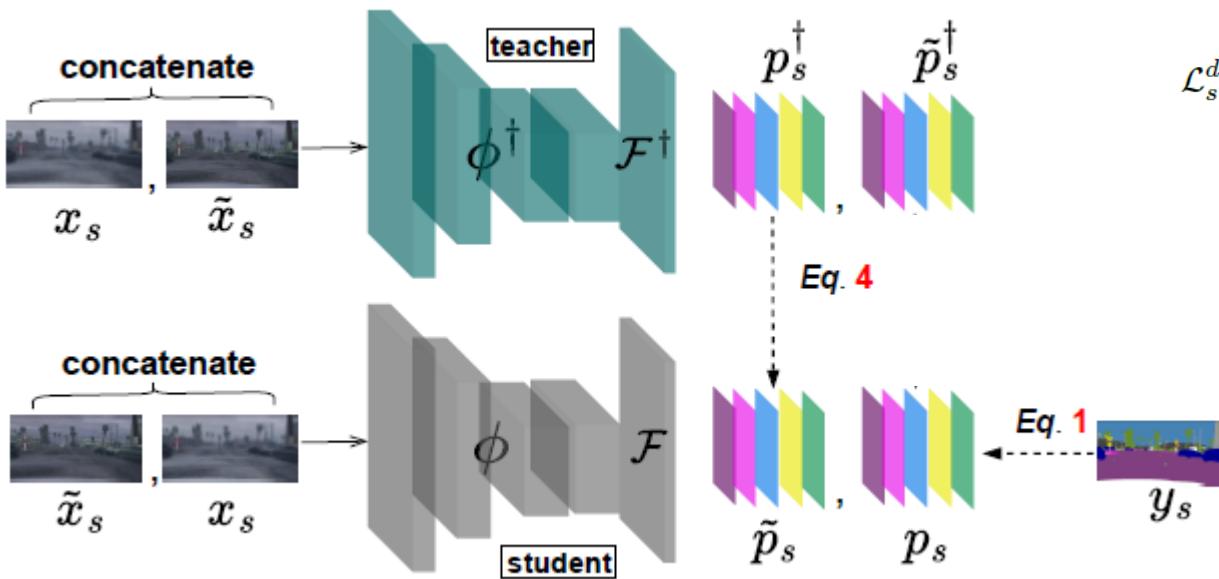
$$\{p_s^\dagger, \tilde{p}_s^\dagger\} = \sigma(\mathcal{F}^\dagger \circ \phi^\dagger(\{x_s, \tilde{x}_s\})) \quad (2)$$

$$\{\tilde{p}_s, p_s\} = \sigma(\mathcal{F} \circ \phi(\{\tilde{x}_s, x_s\})) \quad (3)$$



# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:



$$\mathcal{L}_s^{seg} = \sum_{h,w} \sum_c -y_s^{(c,h,w)} \log(p_s)^{(c,h,w)} \quad (1)$$

$$\{p_s^\dagger, \tilde{p}_s^\dagger\} = \sigma(\mathcal{F}^\dagger \circ \phi^\dagger(\{x_s, \tilde{x}_s\})) \quad (2)$$

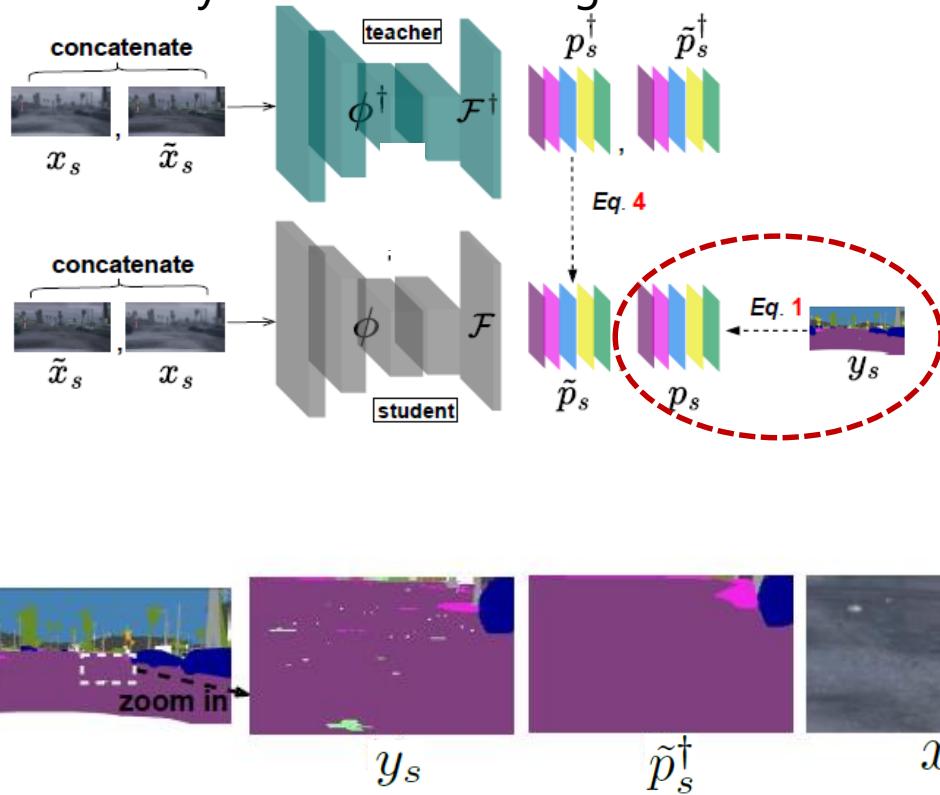
$$\{\tilde{p}_s, p_s\} = \sigma(\mathcal{F} \circ \phi(\{\tilde{x}_s, x_s\})) \quad (3)$$

$$\mathcal{L}_s^{distil} = \overline{\mathcal{H}(p_s^\dagger, \tilde{p}_s)} \quad (4)$$

where  $\mathcal{H}(a,b) = -a \log(b)$

# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:



$$\mathcal{L}_s^{seg} = \sum_{h,w} \sum_c -y_s^{(c,h,w)} \log(p_s)^{(c,h,w)} \quad (1)$$

$$\{p_s^\dagger, \tilde{p}_s^\dagger\} = \sigma(\mathcal{F}^\dagger \circ \phi^\dagger(\{x_s, \tilde{x}_s\})) \quad (2)$$

$$\{\tilde{p}_s, p_s\} = \sigma(\mathcal{F} \circ \phi(\{\tilde{x}_s, x_s\})) \quad (3)$$

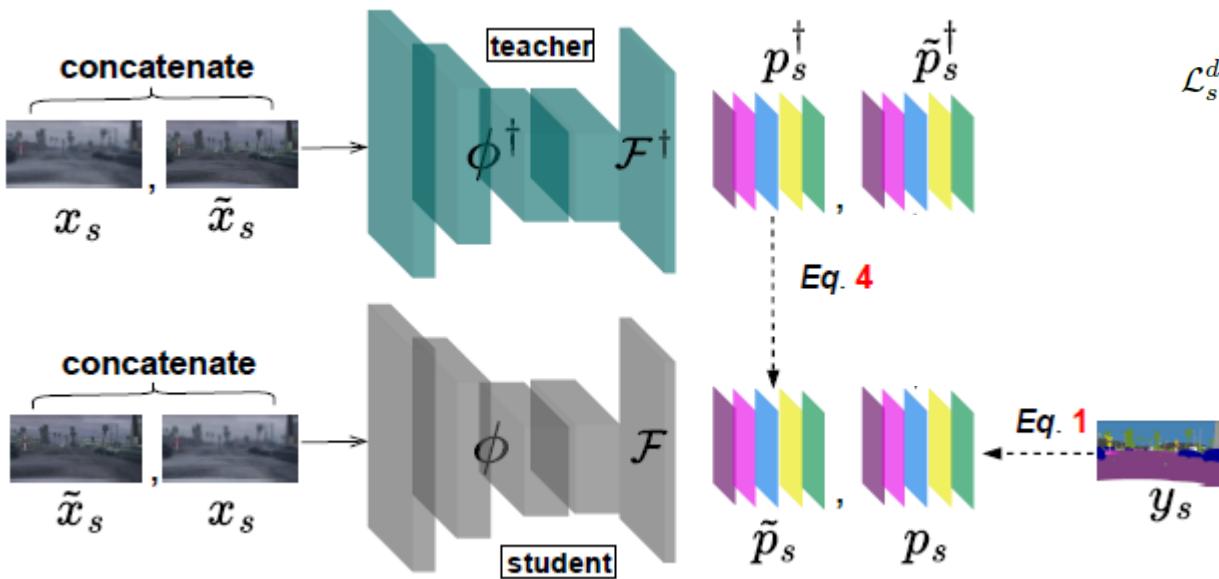
$$\mathcal{L}_s^{distil} = \overline{\mathcal{H}(p_s^\dagger, \tilde{p}_s)} \quad (4)$$

where  $\mathcal{H}(a, b) = -a \log(b)$

Not all source labels are useful for adaptation!

# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:



$$\mathcal{L}_s^{seg} = \sum_{h,w} \sum_c -y_s^{(c,h,w)} \log(p_s)^{(c,h,w)} \quad (1)$$

$$\{p_s^\dagger, \tilde{p}_s^\dagger\} = \sigma(\mathcal{F}^\dagger \circ \phi^\dagger(\{x_s, \tilde{x}_s\})) \quad (2)$$

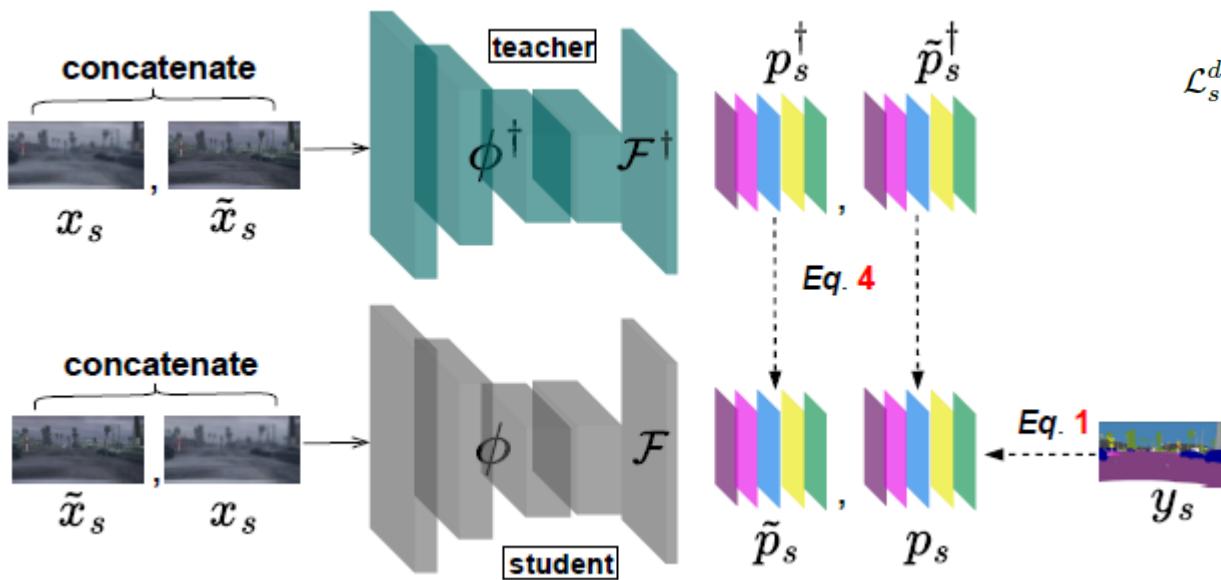
$$\{\tilde{p}_s, p_s\} = \sigma(\mathcal{F} \circ \phi(\{\tilde{x}_s, x_s\})) \quad (3)$$

$$\mathcal{L}_s^{distil} = \overline{\mathcal{H}(p_s^\dagger, \tilde{p}_s)} \quad (4)$$

where  $\mathcal{H}(a,b) = -a \log(b)$

# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:



$$\mathcal{L}_s^{seg} = \sum_{h,w} \sum_c -y_s^{(c,h,w)} \log(p_s)^{(c,h,w)} \quad (1)$$

$$\{p_s^\dagger, \tilde{p}_s^\dagger\} = \sigma(\mathcal{F}^\dagger \circ \phi^\dagger(\{x_s, \tilde{x}_s\})) \quad (2)$$

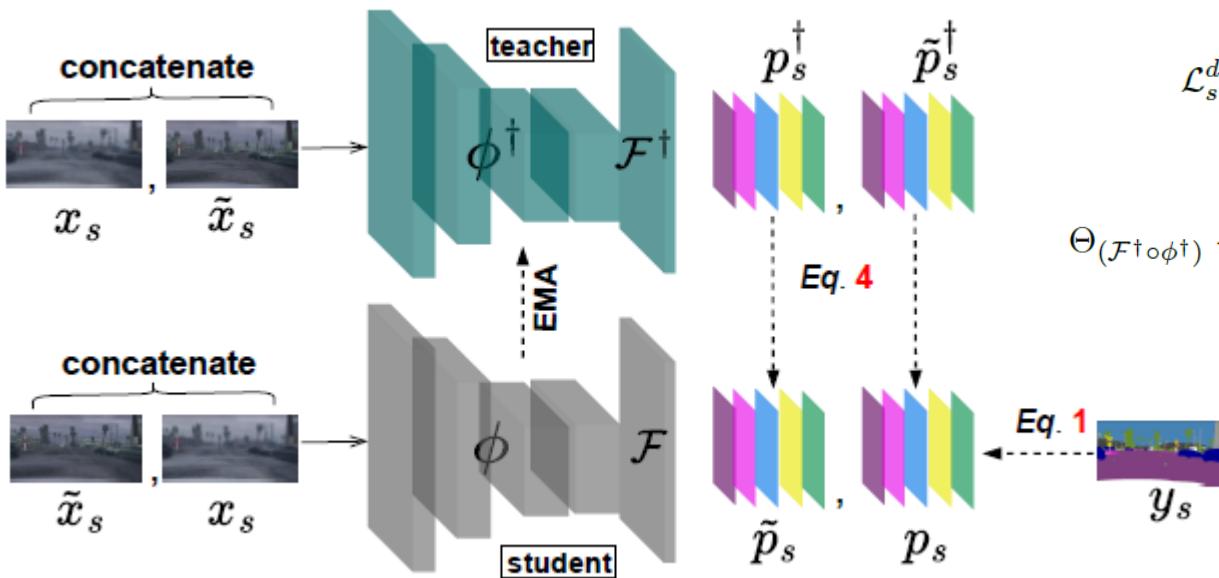
$$\{\tilde{p}_s, p_s\} = \sigma(\mathcal{F} \circ \phi(\{\tilde{x}_s, x_s\})) \quad (3)$$

$$\mathcal{L}_s^{distil} = \overline{\mathcal{H}(p_s^\dagger, \tilde{p}_s^\dagger)} + \alpha \overline{\mathcal{H}(\tilde{p}_s^\dagger, p_s)} \quad (4)$$

where  $\mathcal{H}(a, b) = -a \log(b)$

# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:



$$\mathcal{L}_s^{seg} = \sum_{h,w} \sum_c -y_s^{(c,h,w)} \log(p_s)^{(c,h,w)} \quad (1)$$

$$\{p_s^\dagger, \tilde{p}_s^\dagger\} = \sigma(\mathcal{F}^\dagger \circ \phi^\dagger(\{x_s, \tilde{x}_s\})) \quad (2)$$

$$\{\tilde{p}_s, p_s\} = \sigma(\mathcal{F} \circ \phi(\{\tilde{x}_s, x_s\})) \quad (3)$$

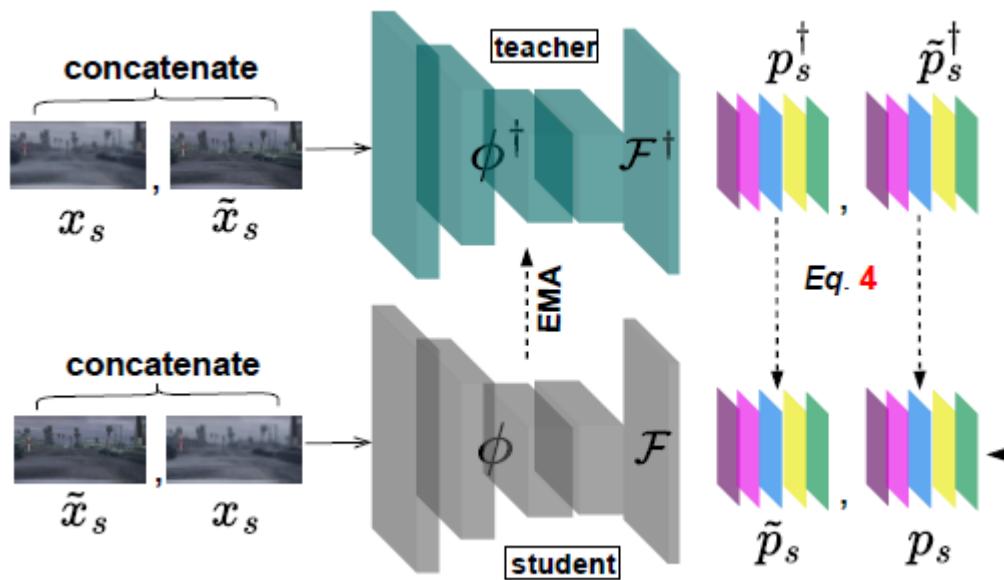
$$\mathcal{L}_s^{distil} = \overline{\mathcal{H}(p_s^\dagger, \tilde{p}_s)} + \alpha \overline{\mathcal{H}(\tilde{p}_s^\dagger, p_s)} \quad (4)$$

$$\text{where } \mathcal{H}(a, b) = -a \log(b)$$

$$\Theta_{(\mathcal{F}^\dagger \circ \phi^\dagger)} \leftarrow \xi * \Theta_{(\mathcal{F}^\dagger \circ \phi^\dagger)} + (1 - \xi) * \Theta_{(\mathcal{F} \circ \phi)} \quad (5)$$

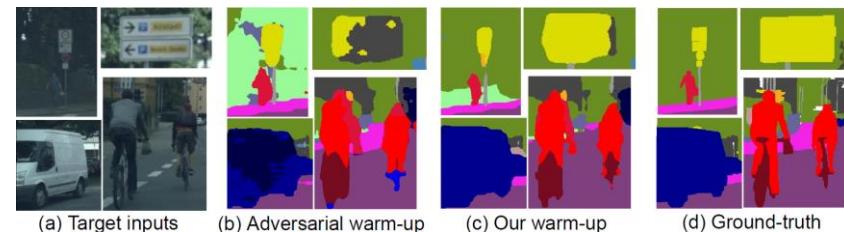
# Method (Warm-up Stage)

Pixel-wise symmetric knowledge distillation:



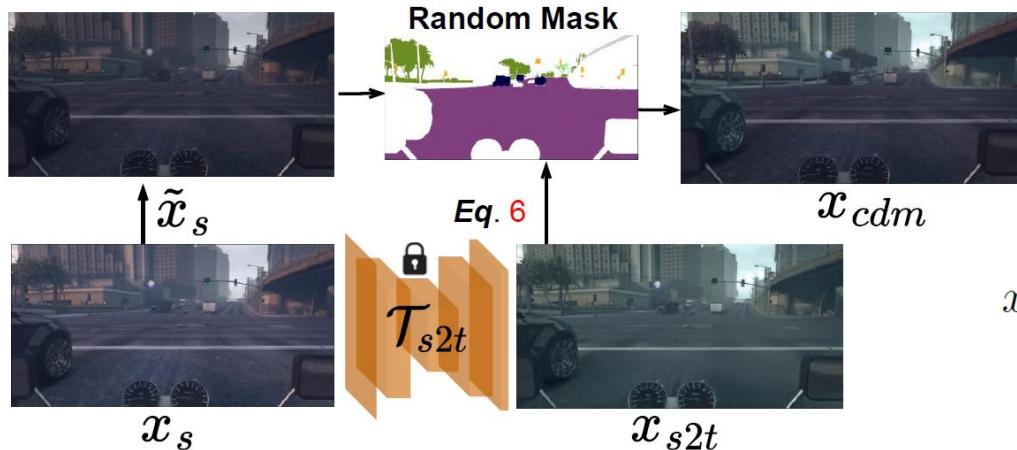
## (+) Benefits are threefold:

- I. knowledge distillation only on source domain, the learning becomes class-aware;
- II. soft labels avoid the model overfitting to domain-specific bias;
- III. our symmetric proposal ensures bidirectional teacher-student consistency between different inputs.



# Method (Warm-up Stage)

## Cross-Domain Mixture Data Augmentation:



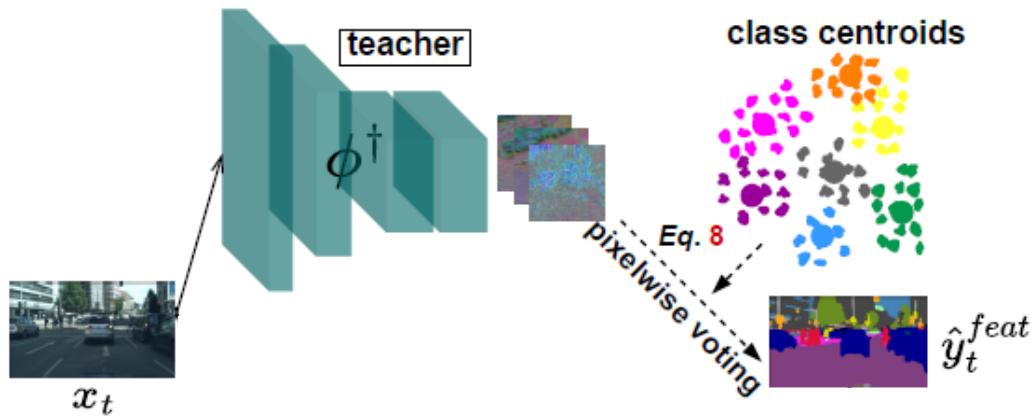
$$x_{cdm} = \tilde{x}_s \odot \mathcal{M} + x_{s2t} \odot (1 - \mathcal{M}) \quad (6)$$

### (+) Benefits :

- I. adding both OOD and target-aware information into data augmentation, distillation more meaningful;
- II. carrying multiple effects on a single view w/o increasing batch size;

# Method (Self-training Stage)

Bilateral-consensus Pseudo-supervision:



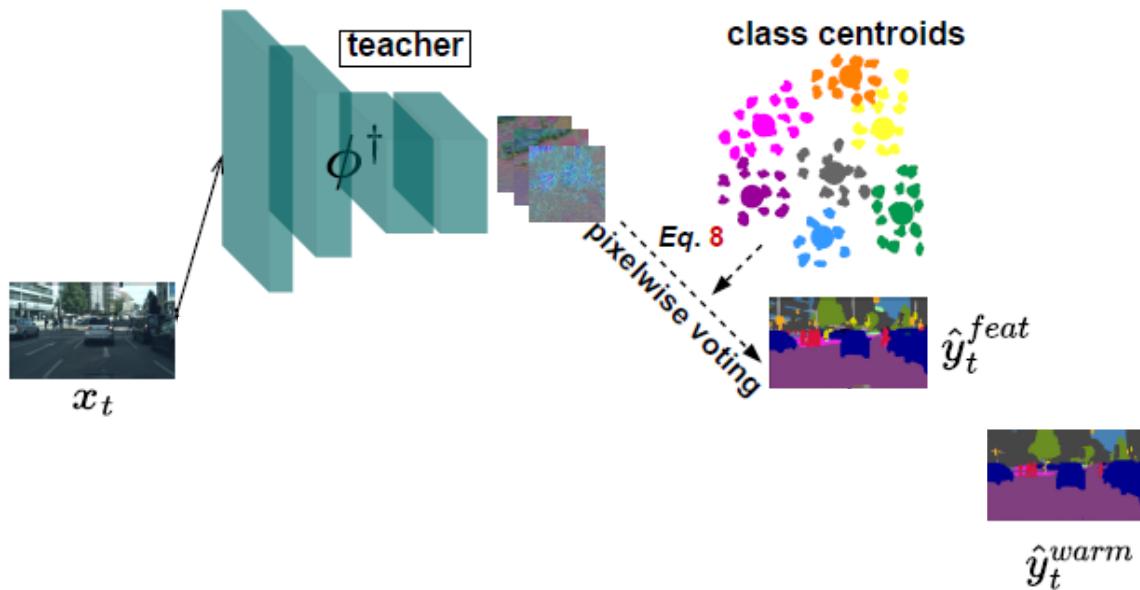
$$\Lambda_s = \{\rho^{(k)}, k = 1, 2, \dots, c\}$$

$$\rho^{(k)} = \frac{\sum_{N_s} GAP(\phi(x_{cdm})^{(k)} \odot (y_s^{(k)} = 1))}{\sum_{N_s} \mathbb{1} \odot (y_s^{(k)} = 1)} \quad (7)$$

$$\hat{y}_t^{feat(jk)} = \mathcal{O}(\arg \min ||\phi^\dagger(x_t)^{(jk)} - \Lambda||_2) \quad (8)$$

# Method (Self-training Stage)

Bilateral-consensus Pseudo-supervision:

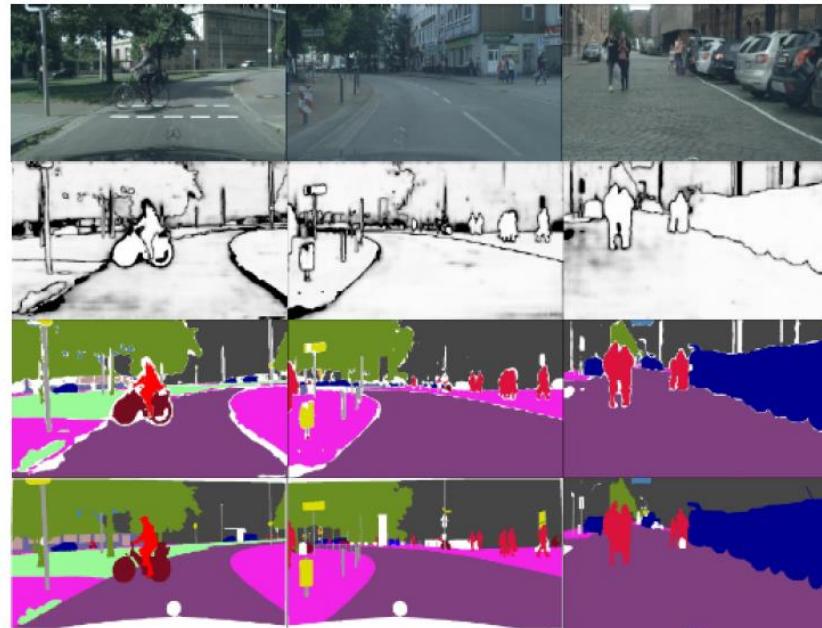
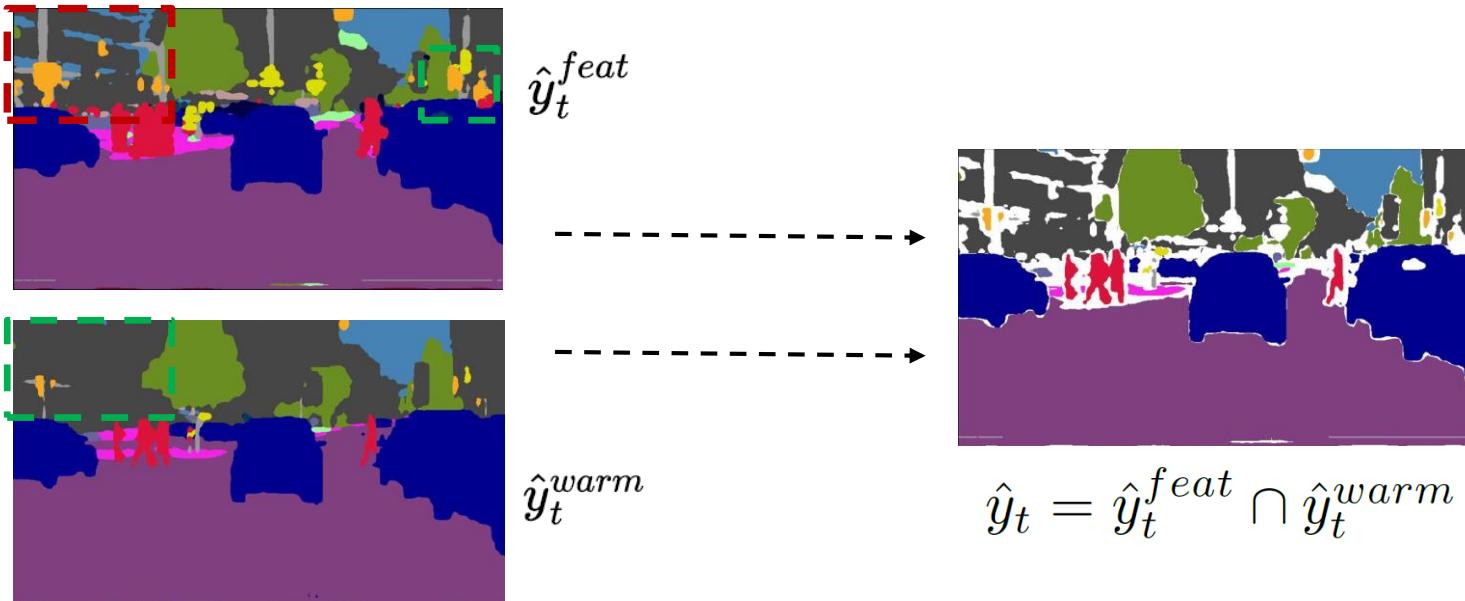


$$\Lambda_s = \{\rho^{(k)}, k = 1, 2, \dots, c\}$$

$$\rho^{(k)} = \frac{\sum_{N_s} GAP(\phi(x_{cdm})^{(k)} \odot (y_s^{(k)} = 1))}{\sum_{N_s} \mathbb{1} \odot (y_s^{(k)} = 1)} \quad (7)$$

$$\hat{y}_t^{feat(jk)} = \mathcal{O}(\arg \min ||\phi^\dagger(x_t)^{(jk)} - \Lambda||_2) \quad (8)$$

## Method (Self-training Stage)



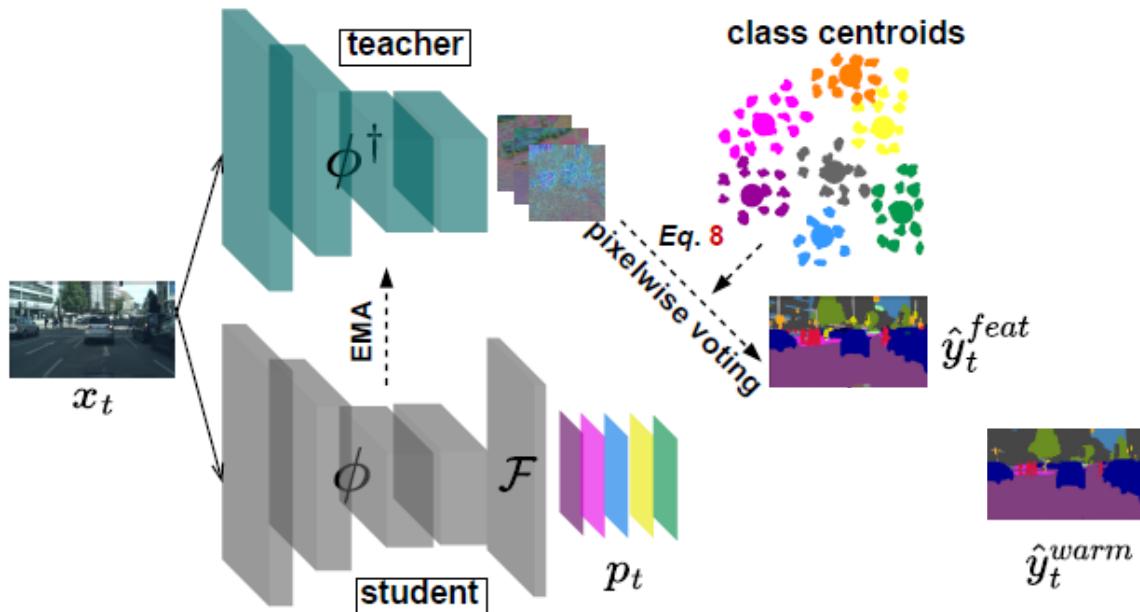
prediction uncertainty  
(darker means lower confidence)

pseudo-label

ground-truth  
(unused in training)

# Method (Self-training Stage)

Bilateral-consensus Pseudo-supervision:



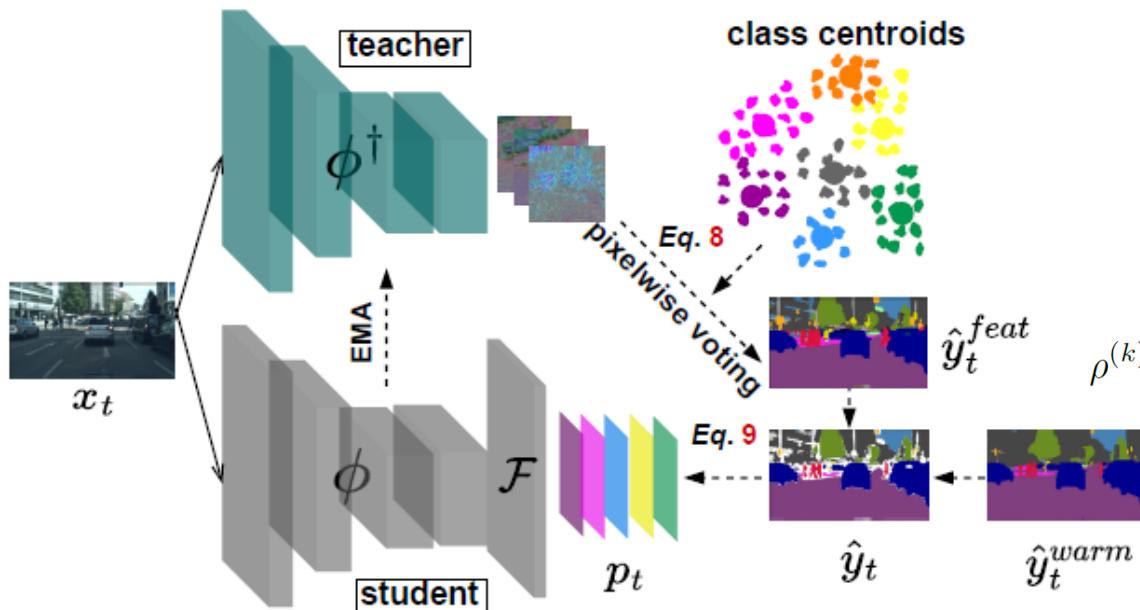
$$\Lambda_s = \{\rho^{(k)}, k = 1, 2..., c\}$$

$$\rho^{(k)} = \frac{\sum_{N_s} GAP(\phi(x_{cdm})^{(k)} \odot (y_s^{(k)} = 1))}{\sum_{N_s} \mathbb{1} \odot (y_s^{(k)} = 1)} \quad (7)$$

$$\hat{y}_t^{feat(jk)} = \mathcal{O}(\arg \min ||\phi^\dagger(x_t)^{(jk)} - \Lambda||_2) \quad (8)$$

# Method (Self-training Stage)

Bilateral-consensus Pseudo-supervision:



$$\Lambda_s = \{\rho^{(k)}, k = 1, 2, \dots, c\}$$

$$\rho^{(k)} = \frac{\sum_{N_s} GAP(\phi(x_{cdm})^{(k)} \odot (y_s^{(k)} = 1))}{\sum_{N_s} \mathbb{1} \odot (y_s^{(k)} = 1)} \quad (7)$$

$$\hat{y}_t^{feat(jk)} = \mathcal{O}(\arg \min \|\phi^\dagger(x_t)^{(jk)} - \Lambda\|_2) \quad (8)$$

$$\hat{\mathcal{L}}_t^{seg} = \sum_{h,w} \sum_c -\hat{y}_t^{(c,h,w)} \log(p_t)^{(c,h,w)} \quad (9)$$

$$\rho^{(k)} \leftarrow \delta(\delta\rho^{(k)} + (1 - \delta)\rho_s'^{(k)}) + (1 - \delta)\rho_t'^{(k)} \quad (10)$$

(+) Benefits :

- I. threshold-free self-training;
- II. dynamic pseudo-label selection;

$$\mathcal{L}_{DiGA} = \lambda_s^{distil} \mathcal{L}_s^{distil} + \lambda^{seg} (\mathcal{L}_s^{seg} + \hat{\mathcal{L}}_t^{seg}) \quad (11)$$

# Experiments (Quantitative Evaluation)

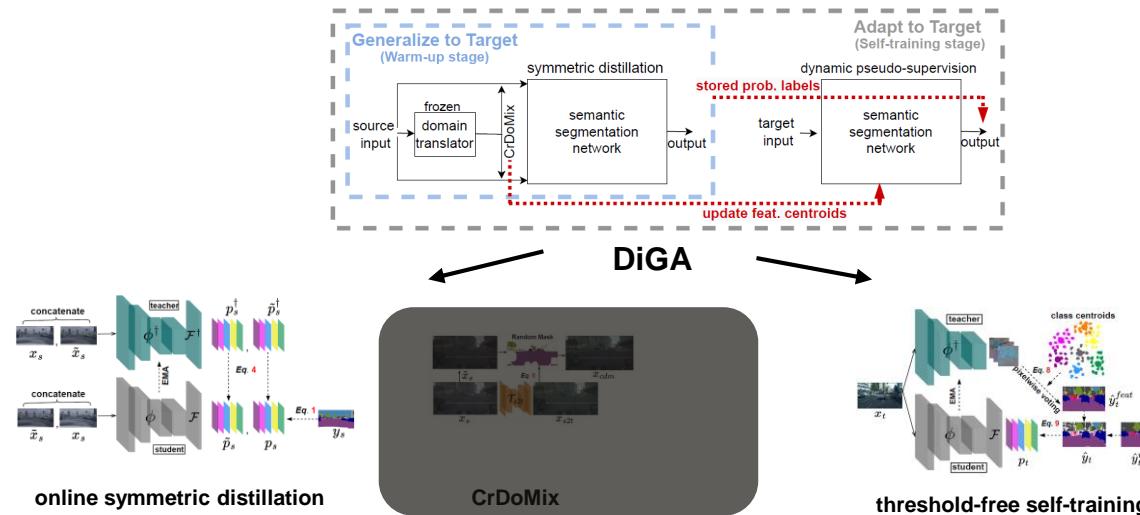
| Method                  | road        | sdkw        | bldng       | wall        | fence       | pole        | light       | sign        | veg         | trn         | sky         | psn         | rider       | car         | truck       | bus         | train       | moto        | bike        | mIoU        |
|-------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| BDL [41]                | 91.0        | 44.7        | 84.2        | 34.6        | 27.6        | 30.2        | 36.0        | 36.0        | 85.0        | <u>43.6</u> | 83.0        | 58.6        | 31.6        | 83.3        | 35.3        | 49.7        | 3.3         | 28.8        | 35.6        | 48.5        |
| ProDA <sup>‡</sup> [80] | 91.5        | 52.4        | 82.9        | 42.0        | <u>35.7</u> | 40.0        | 44.4        | <u>43.3</u> | <b>87.0</b> | <b>43.8</b> | 79.5        | 66.5        | 31.4        | 86.7        | 41.1        | 52.5        | 0.0         | 45.4        | 53.8        | 53.7        |
| CPSL <sup>‡</sup> [40]  | 91.7        | <u>52.9</u> | 83.6        | <u>43.0</u> | 32.3        | <u>43.7</u> | <u>51.3</u> | 42.8        | 85.4        | 37.6        | 81.1        | <u>69.5</u> | 30.0        | 88.1        | <u>44.1</u> | <u>59.9</u> | <u>24.9</u> | <u>47.2</u> | 48.4        | 55.7        |
| ProCA [31]              | 91.9        | 48.4        | 87.3        | 41.5        | 31.8        | 41.9        | 47.9        | 36.7        | 86.5        | 42.3        | 84.7        | 68.4        | 43.1        | 88.1        | 39.6        | 48.8        | <b>40.6</b> | 43.6        | 56.9        | <u>56.3</u> |
| DiGA (Ours, ResNet)     | <b>95.6</b> | <b>67.4</b> | <b>89.8</b> | <b>51.6</b> | <b>38.1</b> | <b>52.0</b> | <b>59.0</b> | <b>51.5</b> | 86.4        | 34.5        | <b>87.7</b> | <b>75.6</b> | <b>48.8</b> | <b>92.5</b> | <b>66.5</b> | <b>63.8</b> | 19.7        | <b>49.6</b> | <b>61.6</b> | <b>62.7</b> |
| DiGA (Ours, HRNet)      | 95.2        | 65.2        | 90.7        | 59.0        | 57.1        | 57.8        | 63.3        | 54.8        | 90.0        | 42.4        | 89.0        | 76.8        | 49.6        | 91.6        | 66.8        | 69.8        | 59.7        | 24.0        | 51.9        | 66.1        |
| DAFormer [27]           | 95.7        | 70.2        | 89.4        | 53.5        | <b>48.1</b> | 49.6        | 55.8        | 59.4        | 89.9        | 47.9        | <b>92.5</b> | 72.2        | 44.7        | 92.3        | 74.5        | 78.2        | 65.1        | 55.9        | 61.8        | 68.3        |
| DiGA (Ours + DAFormer)  | 95.7        | <b>70.4</b> | <b>89.8</b> | <b>54.8</b> | 47.8        | <b>51.3</b> | <b>57.8</b> | <b>63.9</b> | <b>90.3</b> | <b>48.8</b> | 91.8        | <b>73.1</b> | <b>46.6</b> | <b>92.6</b> | <b>78.5</b> | <b>81.3</b> | <b>74.8</b> | <b>57.3</b> | <b>63.2</b> | <b>70.0</b> |
| HRDA [28]               | 96.4        | 74.4        | 91.0        | <b>61.6</b> | 51.5        | <b>57.1</b> | 63.9        | 69.3        | 91.3        | 48.4        | <b>94.2</b> | 79.0        | 52.9        | <b>93.9</b> | <b>84.1</b> | 85.7        | 75.9        | <b>63.9</b> | <b>67.5</b> | 73.8        |
| DiGA (Ours + HRDA)      | <b>97.0</b> | <b>78.6</b> | <b>91.3</b> | 60.8        | <b>56.7</b> | 56.5        | <b>64.4</b> | <b>69.9</b> | <b>91.5</b> | <b>50.8</b> | 93.7        | <b>79.2</b> | <b>55.2</b> | 93.7        | 78.3        | <b>86.9</b> | <b>77.8</b> | 63.7        | 65.8        | <b>74.3</b> |

Table 1. **GTA5-to-Cityscapes adaptation results.** We compare our model performance with state-of-the-art methods. In all tables of Sec. 4.2, bold stands for **best** and underline for second-best. <sup>‡</sup> for fair comparison, we use their reported results after ST stage.

| Method                  | road        | sdkw        | bldng       | wall*       | fence*     | pole*       | light       | sign        | veg         | sky         | psn         | rider       | car         | bus         | mcycl       | bcycl       | mIoU        | mIoU*       |
|-------------------------|-------------|-------------|-------------|-------------|------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| BDL [41]                | 86.0        | 46.7        | 80.3        | -           | -          | -           | 14.1        | 11.6        | 79.2        | 81.3        | 54.1        | 27.9        | 73.7        | 42.2        | 25.7        | 45.3        | -           | 51.4        |
| ProDA <sup>‡</sup> [80] | 87.1        | 44.0        | 83.2        | 26.9        | 0.0        | 42.0        | 45.8        | <u>34.2</u> | 86.7        | 81.3        | 68.4        | 22.1        | 87.7        | 50.0        | 31.4        | 38.6        | 51.9        | 58.5        |
| CPSL <sup>‡</sup> [40]  | 87.3        | 44.4        | 83.8        | 25.0        | 0.4        | <u>42.9</u> | <u>47.5</u> | 32.4        | 86.5        | 83.3        | 69.6        | <u>29.1</u> | <u>89.4</u> | 52.1        | <u>42.6</u> | 54.1        | 54.4        | 61.7        |
| ProCA [31]              | <u>90.5</u> | <u>52.1</u> | 84.6        | <b>29.2</b> | <b>3.3</b> | 40.3        | 37.4        | 27.3        | 86.4        | <b>85.9</b> | <u>69.8</u> | 28.7        | 88.7        | <u>53.7</u> | 14.8        | <u>54.8</u> | 53.0        | 59.6        |
| DiGA (Ours, ResNet)     | 89.1        | <b>53.4</b> | <b>86.1</b> | <u>28.7</u> | <u>3.0</u> | <b>49.6</b> | <b>50.6</b> | <b>34.9</b> | <b>88.2</b> | <u>84.9</u> | <b>71.3</b> | <b>40.9</b> | <b>91.6</b> | <b>75.1</b> | <b>50.3</b> | <b>65.8</b> | <b>60.2</b> | <b>67.9</b> |
| DiGA (Ours, HRNet)      | 90.6        | 56.3        | 87.4        | 38.8        | 6.4        | 57.7        | 59.3        | 50.4        | 87.9        | 86.4        | 76.1        | 47.9        | 89.0        | 54.2        | 47.2        | 69.1        | 62.8        | 69.4        |
| DAFormer [27]           | 84.5        | 40.7        | <b>88.4</b> | 41.5        | 6.5        | 50.0        | 55.0        | <b>54.6</b> | 86.0        | 89.8        | 73.2        | 48.2        | 87.2        | 53.2        | 53.9        | 61.7        | 60.9        | 67.4        |
| DiGA (Ours + DAFormer)  | <b>85.2</b> | <b>41.4</b> | 88.2        | <b>42.6</b> | <b>7.5</b> | <b>52.1</b> | <b>57.5</b> | 47.7        | <b>87.8</b> | <b>90.8</b> | <b>75.0</b> | <b>50.8</b> | <b>87.8</b> | <b>58.0</b> | <b>58.5</b> | <b>63.0</b> | <b>62.1</b> | <b>68.6</b> |
| HRDA [28]               | 85.2        | 47.7        | 88.8        | 49.5        | 4.8        | <b>57.2</b> | <b>65.7</b> | 60.9        | 85.3        | 92.9        | <b>79.4</b> | <b>52.8</b> | 89.0        | <b>64.7</b> | <b>63.9</b> | 64.9        | 65.8        | 72.4        |
| DiGA (Ours + HRDA)      | <b>88.5</b> | <b>49.9</b> | <b>90.1</b> | <b>51.4</b> | <b>6.6</b> | 55.3        | 64.8        | <b>62.7</b> | <b>88.2</b> | <b>93.5</b> | 78.6        | 51.8        | <b>89.5</b> | 62.2        | 61.0        | <b>65.8</b> | <b>66.2</b> | <b>72.8</b> |

Table 2. **Synthia-to-Cityscapes adaptation results.** mIoU, mIoU\* refer to 16-class and 13-class experimental settings, respectively. <sup>‡</sup> for fair comparison, we use their reported results after ST stage following [31].

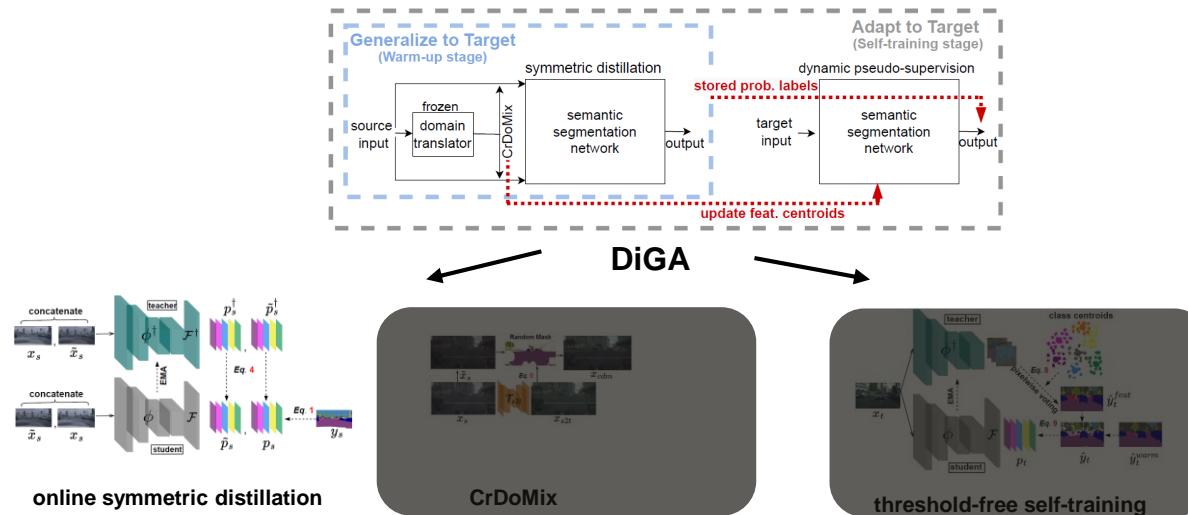
# Experiments (Semi-supervised Semantic Segmentation)



| Method  | Cityscapes   |              |              |              |
|---------|--------------|--------------|--------------|--------------|
|         | 1/16 (186)   | 1/8 (372)    | 1/4 (744)    | 1/2 (1488)   |
| CPS [9] | 75.09        | 77.92        | 79.24        | 80.67        |
| Ours    | <b>76.86</b> | <b>78.51</b> | <b>80.01</b> | <b>80.93</b> |

Table 7. **mIoU comparison of semi-supervised semantic segmentation** using HRNet backbone, based on which SOTA performance of CPS [9] is reported. Evaluation performed on Cityscapes validation set under different partition protocols.

# Experiments (Domain Generalization)



| Method           | Train on GTA5 (G) |                 |                 |                 |
|------------------|-------------------|-----------------|-----------------|-----------------|
|                  | $\rightarrow C$   | $\rightarrow B$ | $\rightarrow M$ | $\rightarrow S$ |
| ISW [12]         | 42.87             | 38.53           | 39.05           | 29.58           |
| SFDA [71]        | 43.50             | -               | -               | -               |
| SAN-SAW [56]     | 45.33             | 41.18           | 40.77           | 31.84           |
| SHADE [83]       | 46.66             | 43.66           | 45.50           | -               |
| Our Distillation | <b>48.87</b>      | <b>44.42</b>    | <b>51.78</b>    | <b>37.17</b>    |

Table 6. **mIoU comparison with SOTA methods for domain generalization.** G, C, B, M and S denote GTA5, Cityscapes, BDD100k, Mapillary and Synthia, respectively. For fair comparison, all the listed methods :

# Experiments (Ablation Study)

## Stage-wise:

| Method      | a | b | c | d | e | mIoU | $\Delta$ |
|-------------|---|---|---|---|---|------|----------|
| Source-only |   |   |   |   |   | 38.3 | +0.0     |
| Source-only | ✓ |   |   |   |   | 38.9 | +0.6     |
| (i)         | ✓ | ✓ |   |   |   | 46.7 | +8.4     |
| (ii)        | ✓ | ✓ | ✓ |   |   | 48.9 | +10.6    |
| (iii)       | ✓ | ✓ | ✓ | ✓ |   | 51.1 | +12.8    |
| (iv)        | ✓ | ✓ | ✓ | ✓ | ✓ | 62.7 | +24.4    |

Table 3. DiGA components: **a**  $\rightarrow$  MST, **b**  $\rightarrow$   $\overline{\mathcal{H}(p_s^\dagger, \tilde{p}_s)}$ , **c**  $\rightarrow$   $\mathcal{H}(\tilde{p}_s^\dagger, p_s)$ , **d**  $\rightarrow$  CrDoMix, and **e**  $\rightarrow \hat{\mathcal{L}}_t^{seg}$ .

## Warm-up stage:

| Strategy | Adv. [67] | Distil. | Adv. [67]+CrDoMix | Distil.+CrDoMix |
|----------|-----------|---------|-------------------|-----------------|
| mIoU     | 45.2      | 48.9    | 47.3              | <b>51.1</b>     |

Table 4. Warm-up model comparison between adversarial training and our knowledge distillation w/ and w/o CrDoMix.

## Self-training stage:

| Strategy | (1) $\hat{y}_t^{feat}$ | (2) $\hat{y}_t^{warm}$ | (3)BDL | (4)ProDA | (5) $\hat{y}_t$ (ours) |
|----------|------------------------|------------------------|--------|----------|------------------------|
| mIoU     | 52.1                   | 53.8                   | 56.2   | 59.5     | <b>62.7</b>            |

Table 5. mIoU comparison of applying different pseudo-labelling techniques to train ST stage **based on our warm-up model**.

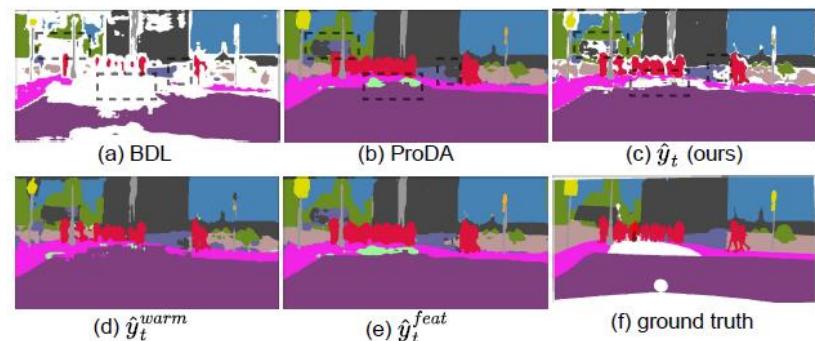


Figure 7. Comparison of different pseudo-labelling techniques given the same input image, and ground truth (f) is only adopted for comparison. Dashed black boxes reveal the major differences.

## Experiments (tSNE, Qualitative)

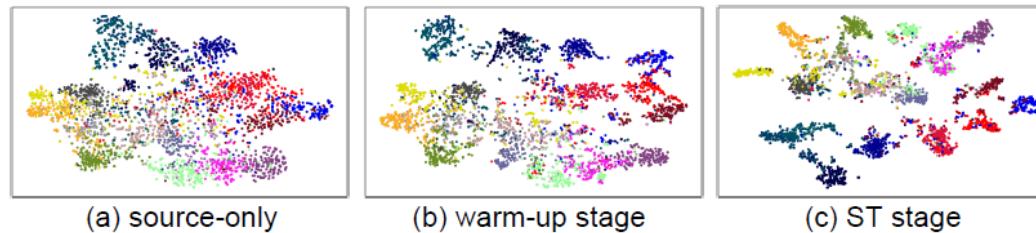


Figure 6. **Visualization of feature distribution** on Cityscapes validation set for each stage based on t-SNE [18] map.

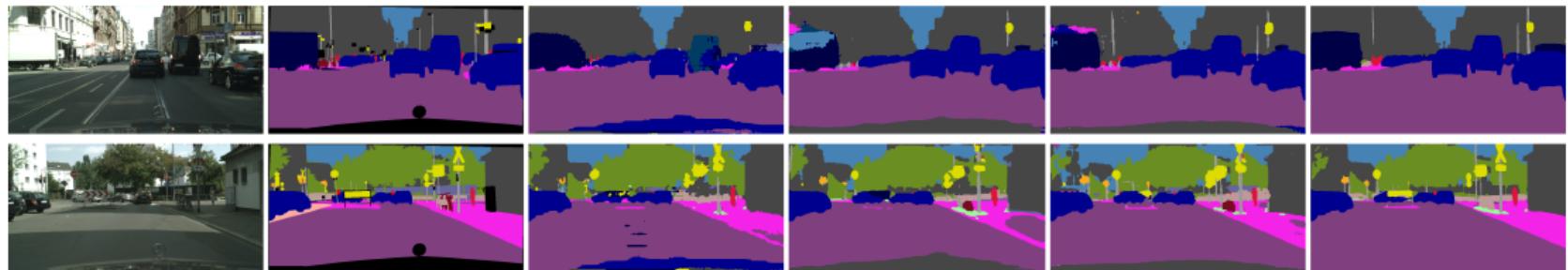


Figure 6. **Qualitative results of GTA5-to-Cityscapes adaptation on Cityscapes validation set.** Columns from left to right are: target domain inputs; ground-truth labels; segmentation predictions of BDL [41], ProDA [80], CPSL [40] and DiGA (ResNet).